

UNITED STATES PATENT APPLICATION
FOR
METHOD OF OBFUSCATING COMPUTER INSTRUCTION STREAMS

INVENTOR:

RAYMOND R. KIDDY

PREPARED BY:

BLAKELY, SOKOLOFF, TAYLOR & ZAFMAN
12400 WILSHIRE BOULEVARD
SEVENTH FLOOR
LOS ANGELES, CA 90025-1026

(408) 720-8300

EXPRESS MAIL CERTIFICATE OF MAILING

"Express Mail" mailing label number: EL 821 772 903 US

Date of Deposit: July 25, 2001

I hereby certify that I am causing this paper or fee to be deposited with the United States Postal Service "Express Mail Post Office to Addressee" service on the date indicated above and that this paper or fee has been addressed to the Assistant Commissioner for Patents, Washington, D. C. 20231

Connie Thayer

(Typed or printed name of person mailing paper or fee)

Connie Thayer

(Signature of person mailing paper or fee)

7-25-01

(Date signed)

METHOD OF OBFUSCATING COMPUTER INSTRUCTION STREAMS

FIELD OF THE INVENTION

The present invention relates to the field of computer science, and more
5 particularly to a method and apparatus for obfuscating computer instruction streams.

BACKGROUND OF THE INVENTION

Some modern compilers, most notably the Java compiler from Sun
Microsystems, are designed to compile source code (e.g. Java Programs or Java
10 Applets) into sequences of instructions to be executed on a stack-based virtual
machine. A key benefit of compiling source code for execution on a virtual machine
is that any processor that can be programmed to implement a virtual machine,
regardless of the processor's internal architecture, may execute the compiled code.

When a human readable unit of source code is compiled into a stream of
15 instructions for a virtual machine by a typical compiler, the mechanically compiled
virtual machine instructions can be deterministically transformed back into a version
of the human readable source code. This process of de-compilation of instructions for
a virtual machine into a version of the human readable source code enables reverse
engineering of the intellectual properties embedded in the source code. After
20 spending a large amount of time and resources developing a software program,
developers do not want to place their applications in the public domain in a form that
gives away their efforts.

Obfuscation is the process of transforming a stream of computer instructions into another stream of instructions that executes the same set of logical operations as that in the original stream so that it is more difficult to be transformed back into a version of the human readable source code.

5 Figure 1 shows one example of an obfuscation method according to one embodiment of the prior art. In operation 341 a typical compiler converts a unit of human readable source code 302 into a virtual machine instruction stream 304 which can be easily de-compiled into a version of the human readable source code. To obfuscate the virtual machine instruction stream 304, operation 343 breaks the stream
10 304 into a set of parts 310. These parts are transformed and padded with dummy instructions in operation 345. For example, part 316 is transformed into part 324, which is padded with dummy instructions 322. The transformations in operation 345 may include reversing loops, expanding loops, flow transformation, renaming identifiers, etc. After the transformation and padding, operation 347 assembles the
15 set of transformed and padded parts 320 into a new instruction stream 330. The new instruction stream is obfuscated and more difficult to be de-compiled into a version of the human readable source code than the mechanically compiled instruction stream 304.

 Dummy instructions 322 are not intended to be executed by a virtual machine
20 for efficiency. For example, null instructions may be used as the dummy instructions to change the patterns of mechanically compiled instruction streams in order to prevent some software programs from de-compiling the instruction stream into a version of the human readable source code.

Figure 2 shows a block diagram of an obfuscation method according to one example of the prior art. Operation 202, corresponding to the operation 343 in Figure 1, breaks a virtual machine instruction stream into parts. Operation 204 transforms the parts; operation 206 pads the transformed parts with dummy instructions.

- 5 Operations 204 and 206 correspond to the operation 345 in Figure 1. Operation 208, corresponding to operation 347 in Figure 1, assembles the padded and transformed parts into a new instruction stream.

- However, the obfuscation methods as in Figures 1 and 2 are subject to attack. The distinct characteristics of the parts, which are taken from a logically cohesive source, and the dummy instructions, which do not perform any logical operation, make it possible to filter out the dummy instructions from the obfuscated instruction stream. Just as chaff can be separated from wheat because of the different physical characteristics, so can be the dummy instructions when an obfuscated instruction stream is compared to an instruction stream that is from a logically cohesive source.
- 10 The chaff can be seen and removed. The dummy instructions may be shown to be garbage or not producible from a valid source, and thus be detected and removed.

- Since in operation 204 the transformations applied to the parts are chosen from a transformation library, a large pool of obfuscated virtual machine instruction streams may be processed to derive the transformation library. With a derived transformation library, an obfuscated instruction stream produced according to the methods in Figures 1 and 2 can be transformed back into a version of a human readable source code once the dummy instructions are removed.
- 20

SUMMARY OF THE INVENTION

Methods and apparatuses for obfuscating computer instruction streams to prevent reverse engineering the human readable source codes of the instruction streams are described here.

5 In one aspect of the present invention, an exemplary method includes breaking each of at least two operative instruction streams into a plurality of parts and interleaving the parts into a new instruction stream.

 In another aspect of the present invention, an exemplary method includes breaking each of at least two operative instruction streams into a plurality of parts and
10 interleaving the parts with obfuscation codes into a new instruction stream. The obfuscation codes interrelate the parts from different instruction streams to prevent reversal of interleaving.

 In another aspect of the present invention, an exemplary method includes breaking each of at least two operative instruction streams into a plurality of parts,
15 transforming and interleaving the parts into a new instruction stream. In one example according to this aspect, the parts are transformed so that the new instruction stream performs at least the same logical operations of one of the operative instruction streams.

 The present invention includes apparatuses which perform these methods,
20 including data processing systems which perform these methods and machine readable media which when executed on data processing system cause the systems to perform these methods.

The present invention also includes machine readable media which contain obfuscated computer instruction streams produced by these methods.

Other features of the present invention will be apparent from the accompanying drawings and from the detailed description which follows.

BRIEF DESCRIPTION OF THE DRAWINGS

The present invention is illustrated by way of example and not limitation in the figures of the accompanying drawings in which like references indicate similar elements.

5 Figure 1 shows a method of obfuscating a computer instruction stream according to one example of the prior art.

Figure 2 shows a block diagram of an obfuscation method according to one example of the prior art.

10 Figure 3 shows a block diagram example of a data processing system which may be used with the present invention.

Figure 4 shows a method of obfuscating computer instruction streams according to one embodiment of the present invention.

Figure 5 shows a block diagram of an obfuscation method according to one embodiment of the present invention.

15 Figure 6 shows another example of obfuscating computer instruction streams according to the present invention.

Figure 7 shows a detailed example of interleaving parts from two computer instruction streams into an obfuscated stream.

20 Figure 8 shows an example of obfuscating computer instruction streams according to the present invention where transformed and interleaved parts from two streams of instructions are interrelated.

Figure 9 shows examples of executing computer instruction streams obfuscated using various methods of the present invention.

Figure 10 shows a block diagram example of executing a combined computer instruction stream.

Figure 11 shows an example of a machine readable media, which may be used to store software and data which when executed by a data processor system causes the
5 system to perform various methods of the present invention.

DETAILED DESCRIPTION

The subject invention will be described with reference to numerous details set forth below, and the accompanying drawings will illustrate the invention. The following description and drawings are illustrative of the invention and are not to be construed as limiting the invention. Numerous specific details are described to provide a thorough understanding of the present invention. However, in certain instances, well known or conventional details are not described in order not to unnecessarily obscure the present invention in detail.

Figure 3 shows one example of a typical computer system which may be used with the present invention. Note that while Figure 3 illustrates various components of a computer system, it is not intended to represent any particular architecture or manner of interconnecting the components as such details are not germane to the present invention. It will also be appreciated that network computers and other data processing systems which have fewer components or perhaps more components may also be used with the present invention. The computer system of Figure 3 may, for example, be an Apple Macintosh computer.

As shown in Figure 3, the computer system 101, which is a form of a data processing system, includes a bus 102 which is coupled to a microprocessor 103 and a ROM 107 and volatile RAM 105 and a non-volatile memory 106. The microprocessor 103, which may be a G3 or G4 microprocessor from Motorola, Inc. or IBM is coupled to cache memory 104 as shown in the example of Figure 3. The bus 102 interconnects these various components together and also interconnects these components 103, 107, 105, and 106 to a display controller and display device 108 and

to peripheral devices such as input/output (I/O) devices which may be mice, keyboards, modems, network interfaces, printers and other devices which are well known in the art. Typically, the input/output devices 110 are coupled to the system through input/output controllers 109. The volatile RAM 105 is typically implemented
5 as dynamic RAM (DRAM) which requires power continually in order to refresh or maintain the data in the memory. The non-volatile memory 106 is typically a magnetic hard drive or a magnetic optical drive or an optical drive or a DVD RAM or other type of memory systems which maintain data even after power is removed from the system. Typically, the non-volatile memory will also be a random access memory
10 although this is not required. While Figure 3 shows that the non-volatile memory is a local device coupled directly to the rest of the components in the data processing system, it will be appreciated that the present invention may utilize a non-volatile memory which is remote from the system, such as a network storage device which is coupled to the data processing system through a network interface such as a modem
15 or Ethernet interface. The bus 102 may include one or more buses connected to each other through various bridges, controllers and/or adapters as is well known in the art. In one embodiment the I/O controller 109 includes a USB (Universal Serial Bus) adapter for controlling USB peripherals.

It will be apparent from this description that aspects of the present invention
20 may be embodied, at least in part, in software. That is, the techniques may be carried out in a computer system or other data processing system in response to its processor, such as a microprocessor, executing sequences of instructions contained in a memory, such as ROM 107, volatile RAM 105, non-volatile memory 106, cache 104 or a

remote storage device. In various embodiments, hardwired circuitry may be used in combination with software instructions to implement the present invention. Thus, the techniques are not limited to any specific combination of hardware circuitry and software nor to any particular source for the instructions executed by the data processing system. In addition, throughout this description, various functions and operations are described as being performed by or caused by software code to simplify description. However, those skilled in the art will recognize what is meant by such expressions is that the functions result from execution of the code by a processor, such as the microprocessor 103.

10 At least one embodiment of the present invention seeks to produce obfuscated streams of virtual machine instructions that are not reversible into a version of human readable source codes. According to the present invention, parts of a number of operative instruction streams are interleaved to produce a combined instruction stream. The combined instruction stream is obfuscated so that it cannot be de-
15 compiled into a version of the human readable source code. Since parts of operative instruction streams are interleaved, each part of the obfuscated stream is a logically cohesive part of an instruction stream. No part of the obfuscated instruction stream can be shown to be garbage or not producible from a valid source. For simplicity, an operative instruction stream will be understood to be a series of instructions
20 corresponding to a method in a source code (e.g., a method in a class, a function, or a subroutine).

 In one embodiment of the present invention, stack-based instruction streams, such as instruction streams for a Java virtual machine, are broken into a number of

stack-balanced blocks. Stack-balanced blocks are instruction sequences that, after execution, leave the stack in the same state as before execution. More specifically, a stack-balanced block is a sequence of instructions that, when executed, causes an equal number of stack push and pop operations to take place in an order such that, throughout execution of the sequence of instructions, the number of completed stack pop operations does not exceed the number of completed stack push operations.

Thus, it is perfectly legal for instructions within a stack-balanced block to push values onto the stack, but the values must be popped off the stack by the end of stack-balanced block so that the stack is unchanged.

When a block of stack-balanced instructions is inserted into an instruction stream, the stack operations caused by the block of stack-balanced instructions do not interfere with the stack operations of the original instruction stream. The local variables used by the block of stack-balanced instructions may be changed so that the operation of the block of stack-balanced instructions has no side effect on the local variables used by the original instruction stream. The local variables and values on the stack may also be transformed within these blocks via functions which, when executed over the course of the entire method in series with the intended function of the method, equal the identity function applied to the intended function of the method. In some embodiments of the present invention, side effects may be intentionally preserved so that the inserted block of stack-balanced instructions enriches the functionality of the original instruction stream. Access to a non-local variable by the block of stack-balanced instructions may be substituted by access to a variable of the same type, which is accessible to the original instruction stream.

Figure 4 shows a method of obfuscating computer instruction streams according to one embodiment of the present invention. In operations 417 and 427 human readable source codes 412 and 422 are mechanically compiled into instruction streams 414 and 424, which can be easily de-compiled into a version of human readable source codes. To obfuscate them, the instruction streams 414 and 424 are broken into parts. The instruction stream 414 is broken into a set of parts 410; the instruction stream 424 is broken into a set of parts 420.

After the instruction streams 414 and 424 are broken into parts, operation 431 interleaves these parts into a new instruction stream 430. For example, parts 426 and 428 taken from the instruction stream 424 are inserted between parts 416 and 418 taken from the instruction stream 414. As a result, in the obfuscated instruction stream 430, parts 446 and 448, which correspond to parts 426 and 428 from the instruction stream 424, are located between parts 436 and 438, which correspond to parts 416 and 418 from the instruction stream 414.

Figure 5 shows a block diagram of an obfuscation method according to one embodiment of the present invention. Operations 512 and 522, corresponding to operations 419 and 429 in Figure 4, break the operative instruction streams into parts. After the instruction streams are broken into parts, the parts are optionally transformed in operations 514 or 524. The optional transformations may involve reversing loops, expanding loops, flow transformation, renaming identifiers, changing the usage of variables, eliminating or substituting instructions, etc. Finally, the optionally transformed parts are interleaved into a new obfuscated instruction stream in operation 532. In other embodiments of the present invention, optional

transformations may also take place before the virtual machine instruction streams are broken into parts.

While Figure 5 shows an example of interleaving two instruction streams into an obfuscated stream, multiple instructions streams can be interleaved into a single
5 obfuscated instructions stream according to the present invention. Figure 6 shows an example where three streams of computer instructions are interleaved.

The order of the parts in an interleaved stream may be different from the order of the corresponding parts in the original stream. Figure 6 shows such an example. Parts from an important stream 610 and unimportant streams 620 and 630 are
10 transformed and interleaved into a new combined stream 640. In Figure 6 it is assumed that part 614 does not depend on part 612. Thus, part 614 can be moved before part 612. In the original instruction stream 610, part 612 is located before part 614. Part 612 is transformed into part 642 in the combined stream 640, and part 614 is transformed into part 644. In the combined stream 640, part 642 is placed after part
15 644. Parts 652 and 662, which are transformed from part 622 of stream 620 and part 632 of stream 630, are inserted between the parts 644 and 642 in the combined stream 640.

In Figure 6 the combined stream 640 contains the parts from the important stream 610, as well as the parts from the unimportant stream 620 and 630. Since the
20 parts from the unimportant streams can also be executed by the virtual machine to perform useful tasks, they don't have to be codes just for the purpose of obfuscation. The unimportant stream 620 and 630 can be compiled from computer programs

which perform certain related tasks. These computer programs may be a part of an application.

Figure 7 shows a detailed example of interleaving parts from two streams into an obfuscated stream. Stream 710, which shows a stream of byte codes for a method of a Java class, is the stream to be obfuscated. Stream 720 is a stream of byte codes for the purpose of obfuscation. Stream 720 may be a stream of byte codes for another method of the same class, or a stream of byte codes for a method of another class, or simply a copy of the stream 710 itself. Stream 710 can be broken into parts 711, 713, 715 and 717. Similarly, stream 720 can be broken into parts 722, 724, 726, and 728. The parts from streams 710 and 720 are interleaved into a stream 730, which performs the same logical operations as the stream 710. To prevent the parts from stream 720 from interfering the operation of the parts from stream 710, a number of transformations are performed. For example, part 722 is transformed into part 732 so that part 732 does not operate on the local variable used by part 731 which is taken from the stream to be obfuscated. Similarly, other parts from streams 710 and 720 are also transformed to avoid interference with each other. If part 728 is placed before part 737, the execution of part 728 makes part 737 not reachable, which is not a desirable side effect. However, if part 728 is placed after part 737, it will not be reachable due to part 737. Therefore, part 728 is discarded.

In one embodiment of the present invention, references to non-local variables in the parts taken from the stream for obfuscation purpose are substituted with references to variables of the same type in the stream to be obfuscated in order to avoid violating any access restriction imposed by a Java Virtual Machine. Calls to

methods that invoke GUI (Graphical User Interface) operations are discarded or substituted with calls to methods that do not create noticeable effects. Some instructions in the parts from the stream for obfuscation purpose may cause control flow changes that may interfere with the proper execution of the stream to be obfuscated. Therefore, these instructions are discarded or substituted with other instructions that have no adverse effects on the proper execution of the stream to be obfuscated. From the above illustration, those skilled in the art can see that various transformations may be applied to the parts to ensure the resulting obfuscated stream is functionally equivalent to the stream to be obfuscated.

Figure 8 shows an example of obfuscating computer instruction streams where transformed and interleaved parts from two streams of instructions are interrelated by obfuscation codes. Obfuscation codes are inserted into the obfuscated stream to relate the parts from different streams to prevent the reversal of interleaving. For example, obfuscation codes 842 and 848 in Figure 8 are inserted into the obfuscated stream 830 to relate the parts from the instruction streams 810 and 820. Parts 812 and 826 are transformed into parts 832 and 836 in the obfuscated stream 830. Obfuscation code 842 is inserted to relate the part 832 from the stream 810 and the part 836 from the stream 820. An obfuscation code may access the variables used by different parts to interrelate them.

Figure 9 shows examples of executing computer streams obfuscated using various methods of the present invention. Figure 9 shows a number of computers, including servers 910, 930, 950 and clients 920, 940, 960. In one scenario, a combined and obfuscated stream, generated according to one of the methods of the

present invention, is transferred from one computer for execution on a virtual machine. For example, the server 950 has mechanically compiled computer instruction streams 951 and 952. The parts of the instruction streams 951 and 952 are interleaved into an obfuscated stream 956. The client 960 downloads the obfuscated stream 956 from server 950 to execute on a virtual machine 967. For instance, server 950 is a web server. The obfuscated stream 956 is a Java application or a Java applet. The client 960 runs a web browser, which downloads the Java application or applet for execution on a virtual machine.

In another scenario, mechanically compiled instruction streams are transferred from a number of computers to a client before a combined and obfuscated stream is produced according to one of the method of the present invention. For example, servers 910 and 930 have instruction streams 911 and 932. The downloaded streams 921 and 922 on client 920 correspond to the streams 911 and 932 on the servers 910 and 930. After the obfuscation process 924, streams 921 and 922 are interleaved into a combined stream, which is executed on a virtual machine 927 on client 920.

In another scenario, mechanically compiled instruction streams, as well as the program which when executed causes a processor to carry to one of the method of the present invention to produce a combined and obfuscated stream, are transferred from a number of computers to a client before the parts from the transferred stream are interleaved into a combined and obfuscated stream. For example, the client 940 downloads an obfuscation program 933 from the server 930, in addition to downloading the streams 911 and 932. Having the downloaded instruction streams 941, 942 and the downloaded obfuscation program 943, client 940 executes the

obfuscation program to carry out obfuscation process 944, which interleaves the parts from the instruction streams 941 and 942 into an obfuscated stream for execution on a virtual machine 947.

Figure 10 shows a block diagram example of executing a combined
5 instruction stream. After receiving from another system a combined stream generated using various methods of the present invention, a computer executes the combined stream. Although Figure 9 or Figure 10 suggests that the client computer receives the obfuscated stream through a network, other media may be used to facilitate the transfer. For example, floppy diskettes, ROM or other removable media may be used
10 to transfer or distribute the combined instruction stream.

Figure 11 shows an example of a machine readable media, which may be used to store software and data which when executed by a data processor system causes the system to perform various methods of the present invention. As noted above, this executable software and data may be stored in various places including for example
15 the ROM 107, the volatile RAM 105, the non-volatile memory 106 and/or the cache 104. Portions of this software and/or data may be stored in any one of these storage devices. The media 1110 for example may be primarily the volatile RAM 105 and the non-volatile memory 106 in one embodiment. The OS 1160 represents an operating system. Instruction streams 1150 and 1140 represent mechanically
20 compiled virtual machine instruction streams. The obfuscated stream 1170 represents the combined stream with parts taken from instruction streams 1150 and 1140. Obfuscation program 1120 represents the computer instructions which when executed by the digital processing system cause the processing system to interleave the parts from operative instruction streams into a combined stream. For example, the parts
25 1152, 1154 and 1156 of the instruction stream 1150 and the parts 1141, 1143 and

1145 of the instruction stream 1140 are interleaved into an obfuscated stream 1170 which has parts 1172, 1174, 1176, 1171, 1173 and 1175. The virtual machine 1130 represents the instructions that implement a virtual machine on the processing system. The combined stream 1170 when executed on the virtual machine 1130 may perform
5 the same set of logical operations as the instruction stream 1150.

In the foregoing specification, the invention has been described with reference to specific exemplary embodiments thereof. It will be evident that various modifications may be made thereto without departing from the broader spirit and scope of the invention as set forth in the following claims. The specification and
10 drawings are, accordingly, to be regarded in an illustrative sense rather than a restrictive sense.